



ψ沈黙のジャスティスψ

PWS Cup 2015

全部門第一位チーム
「ψ沈黙のジャスティスψ」が解説する
最強の匿名加工技術

NTTセキュアプラットフォーム研究所
濱田 浩気

概要



Innovative R&D by NTT

昨年の PWS Cup で、ψ沈黙のジャスティスψは
幸運にもたくさんの賞をいただきました。
その中で何をしていたのかをご紹介します。



- ψ 沈黙のジャスティス ψ (概要, 方針)
- 匿名加工フェーズ (予備戦, 本戦)
- 再識別フェーズ (予備戦, 本戦)

ψ沈黙のジャスティスψとは



- NTT という会社の真面目な3人の社員



- 表の顔: プライバシー保護技術の研究者
- 昨年7月～: ψ沈黙のジャスティスψで大会参加

- (かなり)勝負に徹する
 - 会社からのプレッシャー=応援
 - やりすぎと思われた方, すみません...
- でも大会は盛り上げたいし楽しみたい
 - 匿名加工は各自一つずつ好きなものを提出
 - 再識別フェーズを考えると, 本命以外の2枠は当てられにくいものを提出するのが有利だが
 - ヒントは与える(後述)が, 早めにいくつか提出

- 相対評価

- どれを提出していいかわからない

- 匿名加工では多様なデータを用意・提出し、良いものを残そう

- 評価値は公開

- 早めの提出は他チームへの手法のヒントに

- できるだけ遅めに提出しよう、
他チームの評価値をよく見よう

- サンプルの再識別率, U5も公開
 - 匿名加工データからは分かり得ない情報
実世界では得られないが, 再識別に有用
 - 評価値を積極的に利用
- あくまで指標値で決着
 - 既存の手法に囚われず,
指標値の最適化に集中

匿名加工フェーズ(予備戦)

おおまかな戦術の選定



• いくつかの指標で勝てそうかで見積もり

– 対象: 元のまま, k匿名, YA, etc.

	重み	指標	元のまま	k匿名	YA
有用性	2	U1(meanMAE)	◎	○	◎
	2	U2(crossMean)	◎	△	◎
	2	U3(crossCount)	◎	△	◎
	2	U4(corMAE)	◎	×	◎
	2	U5(IL)	◎	×	△
	2	U6(nrow)	◎	◎	◎
安全性	3	S1(min-k)	△	◎	△
	3	S2(avg-k)	△	◎	△
	6	S3(max-reid)	×	○	◎
評価			14	14.7	18.7

凡例

◎: 1位

○: 上から1/3

△: 上から2/3

×: 最下位近く

具体例1: フリーダム(予備戦1位)



申請番号	チーム	提出ファイル	A)有用性U1~U6 平均ランク	B)安全性S1~S2 ランク合計/4	C)安全性S3~S6 最高値ランク/2	総合点 A+B+C
1513	ψ沈黙のジャスティスψ	[フリーダム]	3.00000	2.75000	0.50000	6.25000
1494	SSTK	[SSTK_ICE_FINAL04]	3.16667	2.75000	0.50000	6.41667
1515	MDL	[演员的自我修养]	3.66667	2.75000	0.50000	6.91667
1490	SSTK	[SSTK_ICE_FINAL03]	3.83333	2.75000	0.50000	7.08333
1528	MDL	[ちゃんとしてない]	4.00000	2.75000	0.50000	7.25000
1426	Tsukuba-KDE	[CCM02]	5.00000	2.75000	0.50000	8.25000
1505	Tsukuba-KDE	[CCM03]	5.00000	2.75000	4.00000	11.75000
1424	チームすててこ伊藤	[ステテコY3]	6.00000	8.75000	0.50000	15.25000
1583	ψ沈黙のジャスティスψ	[五十万郎丸]	4.66667	2.75000	8.00000	15.41667
1152	@kusano_k	[RhoAias2]	4.83333	2.75000	8.00000	15.58333
1589	ψ沈黙のジャスティスψ	[ジャスティスなピーパーの大工事]	11.16667	2.75000	4.00000	17.91667
1151	@kusano_k	[RhoAias1]	1.00000	2.75000	15.50000	19.25000
1493	Tsukuba-KDE	[わらび*もちもち]	13.66667	2.75000	4.00000	20.41667
1496	SSTK	[SSTK_ICE_FINAL06]	13.16667	2.75000	6.50000	22.41667
1154	@kusano_k	[RhoAias3]	1.00000	8.75000	15.50000	25.25000
1518	nifigaki	[hoge]	1.00000	8.75000	15.50000	25.25000
1519	nifigaki	[a1]	1.00000	8.75000	15.50000	25.25000
1568	圧倒的「成長」	[ファイナルエスカレーション5]	18.16667	2.00000	6.50000	26.66667
1498	MDL	[ちゃんとしてる]	12.50000	2.75000	11.50000	26.75000
1473	圧倒的「成長」	[リバイズ5]	15.00000	2.75000	9.50000	27.25000

フリーダム(予備戦1位)でしたこと



- QI統一
- YA (+IL最適化(+ランダム化))

	重み	指標	YA		フリーダム
有用性	2	U1(meanMAE)	◎		◎
	2	U2(crossMean)	◎		◎
	2	U3(crossCount)	◎		◎
	2	U4(corMAE)	◎		◎
	2	U5(IL)	△	→	○
	2	U6(nrow)	◎		◎
安全性	3	S1(min-k)	△		△
	3	S2(avg-k)	△	→	○
	6	S3(max-reid)	◎		◎
評価			18.7	→	20.3

凡例

◎: 1位

○: 上から1/3

△: 上から2/3

×: 最下位近く

- 手法: U2, U3 で 対象外のカテゴリ属性 の値を 単一 にする
- 効果: ほぼノーリスクで S2 を高められる

世帯	性別	産業	職業	食料支出	
1	1	2	5	29496.0792	...
1	2	2	5	25806.1846	...
1	2	3	6	38278.1598	...
2	2	3	VV	74122.0521	...
2	1	1	5	33256.8355	...
2	2	1	6	46992.7870	...
⋮	⋮	⋮	⋮	⋮	



世帯	性別	産業	職業	食料支出	
1	1	2	5	29496.0792	...
1	1	2	5	25806.1846	...
1	1	3	5	38278.1598	...
2	1	3	5	74122.0521	...
2	1	1	5	33256.8355	...
2	1	1	5	46992.7870	...
⋮	⋮	⋮	⋮	⋮	

- 手法: 行ごとに属性値を入れ替え
- 効果: U5 は失うが, U1-U4,U6,S1,S2 を維持しながら S3 を高められる

行番号	性別	産業	職業	食料支出	
1	1	2	5	29496.0792	...
2	2	2	5	25806.1846	...
3	2	3	6	38278.1598	...
4	2	3	VV	74122.0521	...
5	1	1	5	33256.8355	...
6	2	1	6	46992.7870	...
⋮	⋮	⋮	⋮	⋮	



行番号	性別	産業	職業	食料支出	
1	1	1	5	33256.8355	...
2	2	2	5	25806.1846	...
3	2	3	6	38278.1598	...
4	2	1	6	46992.7870	...
5	2	3	VV	74122.0521	...
6	1	2	5	29496.0792	...
⋮	⋮	⋮	⋮	⋮	

- 手法: YA の高度化. U5 が大きくなならない ように行ごとに属性値を入れ替え
- 効果: YA で U5 も小さめ.
- 実装方針: 匿名化対象データは8333レコードしかないないので, $O(n^2)$ でもなんとかなるはず
 - 各行間の距離(IL)をすべて計算し,
小さいペアから採用する $O(n^2 \log n)$ の貪欲法
– 簡単な枝刈りをしつつで1分程度

- **手法:** 「YA + IL最適化」の高度化
 - 貪欲法をランダム化し, 乱数の範囲を調整して U5 と S3 のバランスを取れるようにする
- **効果:** U5 は少し大きくなるが, 当てられにくくなる
- **実装:** 行間の距離に ノイズ(乱数)を加算
 - ノイズ大 → 当てられにくく, U5 大きく
 - ノイズ小 → 当てられやすく, U5 小さく

匿名加工フェーズ(本戦)

具体例2: アイロン(本戦1位)



事務局へのデータ提出リスト (総合点数)

チーム	提出ファイル	A) $U_1 \sim U_6$	B) $S_1 \sim S_2$	C) $E_1 \sim E_5 + \alpha$	総合点 A+B+C
ψ沈黙のジャスティスψ	chinmoku_1	1.33333	1.25000	0.50000	3.08333
ψ沈黙のジャスティスψ	chinmoku_3	1.50000	1.25000	1.00000	3.75000
ψ沈黙のジャスティスψ	chinmoku_2	1.66667	1.25000	1.50000	4.41667
ψ沈黙のジャスティスψ	chinmoku_4	3.00000	1.00000	2.00000	6.00000
nifigaki	nifigaki_1	3.16667	1.75000	2.50000	7.41667

2015/10/21

アイロン(本戦1位)でしたこと



- フリーダム(QI統一, YA+IL最適化+ランダム化)
- QIグループ内スワップ

	重み	指標	フリーダム		アイロン
有用性	2	U1(meanMAE)	◎		◎
	2	U2(crossMean)	◎		◎
	2	U3(crossCount)	◎		◎
	2	U4(corMAE)	◎	→	×
	2	U5(IL)	○		○
	2	U6(nrow)	◎		◎
安全性	3	S1(min-k)	△		△
	3	S2(avg-k)	○		○
	6	S3(max-reid)	(◎→)×	→	◎
評価			(20.3→)14.3	→	18.3

Q1グループ内スワップ(1/2)



- 手法: 各Q1グループ内で数値属性値を
属性ごとに独立にスワップ

世帯	人員	消費支出	食料	住居
1	1	2398.98	193.23	345.45
1	1	2389.11	78.22	567.88
1	1	2092.98	987.64	412.23
2	1	2888.10	234.23	452.90
2	1	2074.42	105.32	576.12
2	1	2615.21	688.98	662.21
⋮		⋮	⋮	⋮



世帯	人員	消費支出	食料	住居
1	1	2092.98	78.22	412.23
1	1	2398.98	987.64	345.45
1	1	2389.11	193.23	567.88
2	1	2074.42	234.23	576.12
2	1	2615.21	688.98	452.90
2	1	2888.10	105.32	662.21
⋮		⋮	⋮	⋮

• 効果:

- AYA に強くなった(S3 向上)
- 同一Q1グループ内でのスワップなので...
 - U2, U3 には影響なし
 - U4 は悪化するが諦める

• なぜうまくいくのか:

- AYA は総和に基づくソートで推定
 - 総和が乱れると隙ができる(AYAの推定が甘くなる)
 - その後IL最適化を行えば, この隙に入り込める

- 数値属性で分けたクラス中での入れ替え
(五十万郎丸(予備戦9位))
 - その際に数値を平均化 = 総和を乱して AYA 対策
(Anony Name(本戦2位?), ponnyo(本戦4位?))
- 平均と相関を保存する擬似データを生成
(ジャスティスなビーバーの大工事(予備戦11位))
- Q1グループ間のレコード移動で S2(min-k) 向上
 - U2 をあまり変えないレコードを選んで移動
 - U3 を壊すので思ったほど効果が上がらず

再識別フェーズ

データを眺めて3つに分類

1. 微小なノイズとかななどで簡単に解けそう
→ 簡単に解く
2. サンプルで十分によさそう
→ サンプルを適用
3. YA みたいだから無理そう
→ IL 近傍推定
→ 教科書 YA 解読

- 発想: YA を使う場合は IL (U5) を 小さく するはず
→ IL が小さいレコード を選んで推定結果とすれば、
そこそこ当たるのでは
- 手法: 各レコードに対し、ILが最小の行番号 を推定値として出力
- 実装: 8333行なので $O(n^2)$ の実装. 6秒くらい

経緯:

- YA 採用チームはきっとU5を小さくしているはずなのに U5 が 0.037 と大きいデータがある？
- ランダムでも U5=0.035 程度のはずなのに??
- もしかして、教科書(ルール説明論文)
そのままの実装があったりするのでは？
→ あ、あった！(U5=0.028だけど)
- さらに、ずれ幅を変えていたりとか？
→ あ、あった！(U5=0.037)

• 手法:

ずれ幅1~8332の教科書YAすべてのILを計算
一致する評価値を持つ匿名加工データを探す

頑張って計算したIL

ずれ	IL (U5)
1	0.02796416
2	0.02996592
3	0.03217687
4	0.03292660
5	0.03330663
6	0.03331962
⋮	⋮
8332	0.02796416

見~つけた!

各匿名加工データの指標値

名前	U1	...	U5	...
データ1	0.00		0.01305021	...
データ2	1.01		0.00000000	...
データ3	0.00		0.02796416	...
⋮	⋮		⋮	

ψ沈黙のジャスティスψのしたことを解説

指標値の最適化に集中した結果:

- 匿名加工:
 - YA をベースに IL最適化
 - Q1グループ内スワップ
- 再識別:
 - IL近傍推定
 - 教科書YA解読

- とても勉強になりましたし、
何より、非常に楽しいコンテストでした
- 一緒に競い合ってくださったみなさま、
運営のみなさま、誠にありがとうございました